

# Data Wrangling

The data for this project was gathered from different sources, which are the

- Twitter API
- Udacity website.

The data consist of three different datasets: ' df\_twitter table, ` df\_dog\_image table, ` and ` df\_dog\_ table.` Each dataset contains additional data which are related to each other.

## 1. Twitter API

Access to the Twitter API privileges is given to a Twitter developer account, where access to users' tweets and favorite count, and retweet count can be collected and analyzed by a third party. The data is given in JSON format.

## 2. Udacity website

- Twitter\_archive.csv:

The Udacity website provides access to the tweets from the @welovedogs Twitter account from their archives. The WeRateDogs Twitter archive contains primary data for all 5000+ of their tweets. This information is contained in a CSV file called Twitter\_archive.csv.

- Image\_prediction.tsv :

Images posted on the @welovedogs Twitter account were trained by running through a neural network to classify the breeds of dogs. The result is a table full of the top three image predictions alongside each tweet ID, image URL, and the image number corresponding to the most confident predictions.

**Df\_dog** Table

• <b>Tweet_id</b>	unique The id number of a tweet
• <b>In_reply_to_status_id</b>	The unique id number of a tweet reply
• <b>in_reply_to_user_id</b>	The unique id of the user replied
• <b>timestamp</b>	The timestamp of the tweet
• <b>Source</b>	The source of the tweet
• <b>text</b>	The tweet contents
• <b>retweeted_status_id</b>	Unique tweet for status retweet
• <b>retweeted_status_user_id</b>	Unique user id for status retweet
• <b>expanded_urls</b>	links
• <b>rating_numerator</b>	ratings
• <b>name</b>	Dogs name
• <b>doggo</b>	
• <b>pupper</b>	
• <b>floofer</b>	
• <b>puppo</b>	

**Df\_Twitter** Table

Tweet_id	Unique_tweet id
date_created	Date Tweeted
retweet_count	Number of retweets
Favourite_count	Number of favorites

**Df\_dog\_image** Table

<b>tweet_id</b>	Unique tweet id
<b>jpg_url</b>	Link to image
<b>Image number</b>	The image number predicted as a tweet can have more than one image(up to 4)
<b>p1</b>	P1 is the algorithm's #1 prediction for the image in the tweet
<b>p1_conf</b>	p1_conf is how confident the algorithm is in its #1 prediction
<b>p1_dog</b>	p1_dog is whether or not the #1 prediction is a breed of dog
<b>p2</b>	P2 is the algorithm's #2 prediction for the image in the tweet
<b>p2_conf</b>	p2_conf is how confident the algorithm is in its #2 prediction
<b>p2_dog</b>	p2_dog is whether or not the #2 prediction is a breed of dog
<b>p3</b>	P3 is the algorithm's #3 prediction for the image in the tweet
<b>p3_conf</b>	p3_conf is how confident the algorithm is in its #3 prediction
<b>p3_dog</b>	p3_dog is whether or not the #3 prediction is a breed of dog